# Big Data Management for Grids and Clouds:
## *Lessons Learned by Particle Physics to Take to the Future*

**Adam Lyon / Fermilab SCD / Muon g-2**
**ISGC, Taipei, March 2014**

# Some context

**Experimental Particle Physics in 2 minutes...**
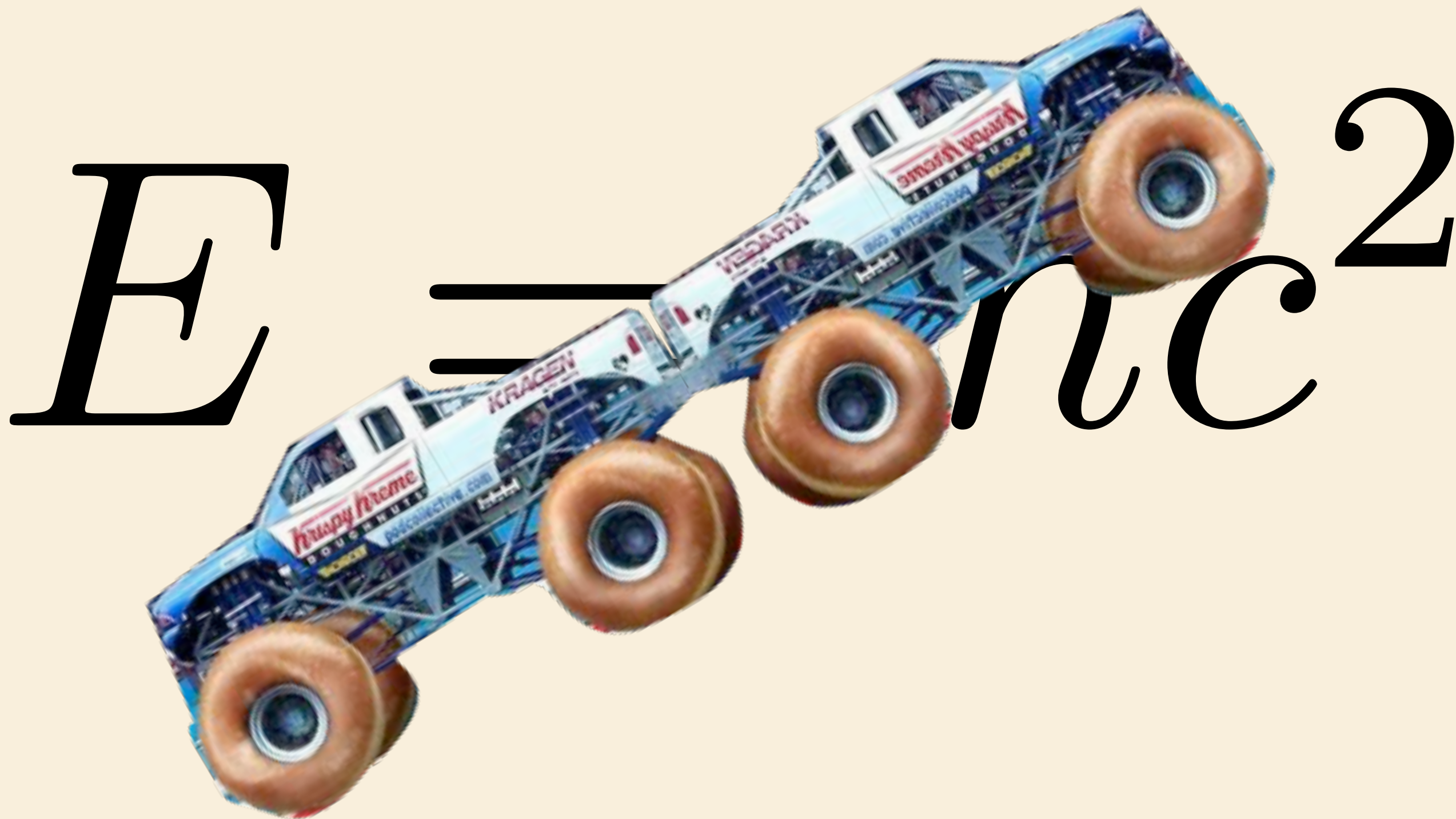
Fermilab

# Collisions

$$E = mc^2$$

❖ **Fermilab**

# Collisions

🍴 Fermilab

# Collisions

$$E = mc^2$$

Data Management Lessons from Particle Physics | Adam Lyon

2014 March 28

**Fermilab**

# Collisions

**Fermilab**

# The Energy Frontier



CDF

Fermilab
Tevatron
4 mi, 2 TeV
1983-2011

D0



ALICE

CMS

ATLAS

LHCb

CERN
Large Hadron Collider
17 mi, 7–14 TeV
2010–

🎱 Fermilab

# Neutrino Oscillations

Fermilab

# Neutrino Oscillations

**✦ Fermilab**

# Neutrino Oscillations

NOvA, MicroBoone, LBNE
Dya Bay, T2K, OPERA

**🎇 Fermilab**

# Muon *g-2*

**Muon in a magnetic storage ring**

➡ **Muon travel direction**

→ **Muon spin direction**

We don't see this

Instead, we see this

# Muon *g-2*

**Muon in a magnetic storage ring**

→ **Muon travel direction**

→ **Muon spin direction**

We don't see this

Instead, we see this

🎗 **Fermilab**

# Muon *g-2*

**Muon in a magnetic storage ring**

→ **Muon travel direction**

→ **Muon spin direction**

We don't see this

Instead, we see this

Data Management Lessons from Particle Physics | Adam Lyon
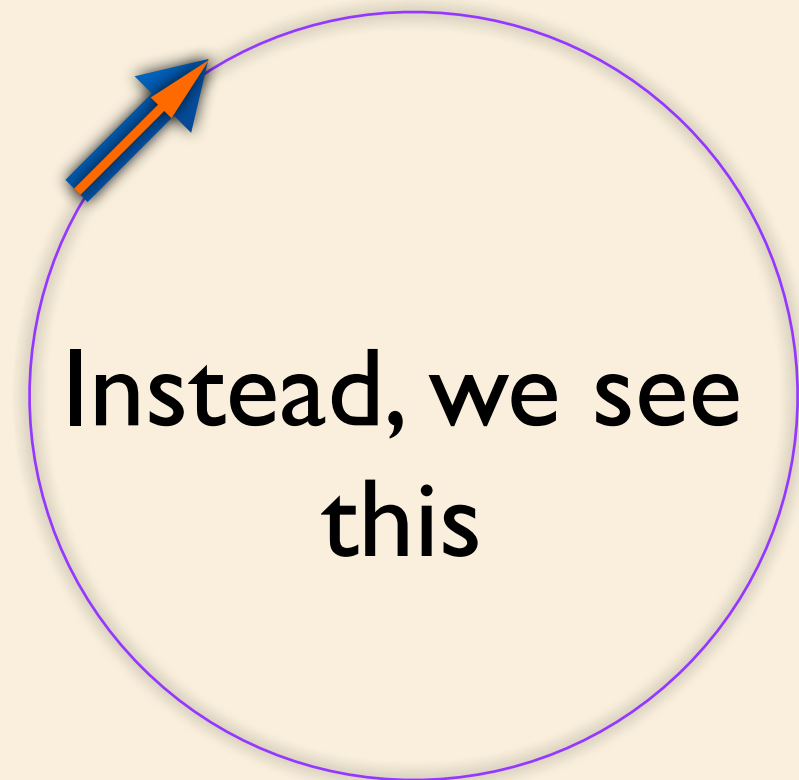
2014 March 28

🌼 **Fermilab**
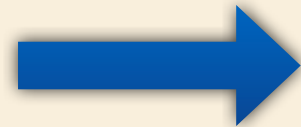
# Muon *g-2*
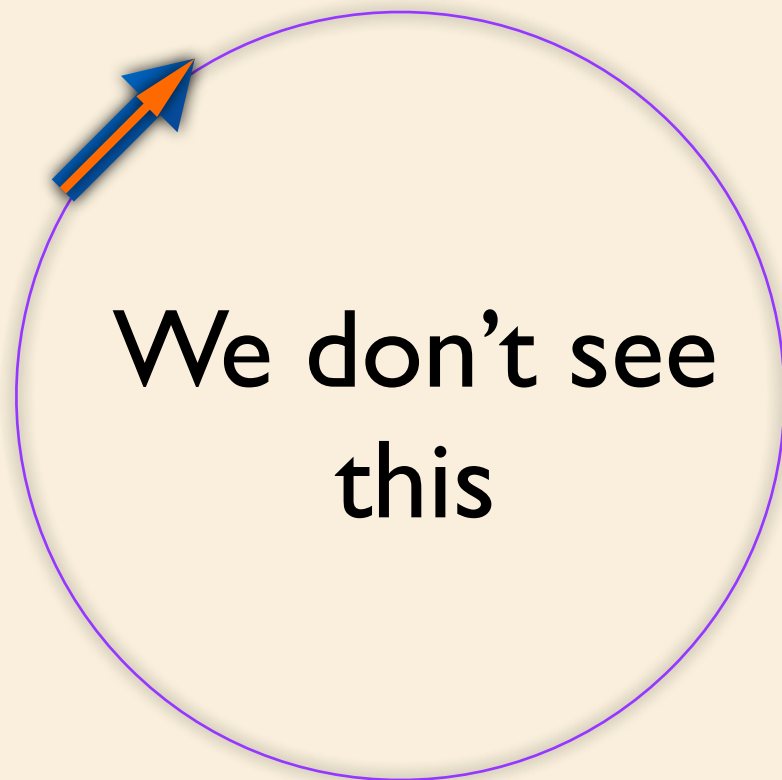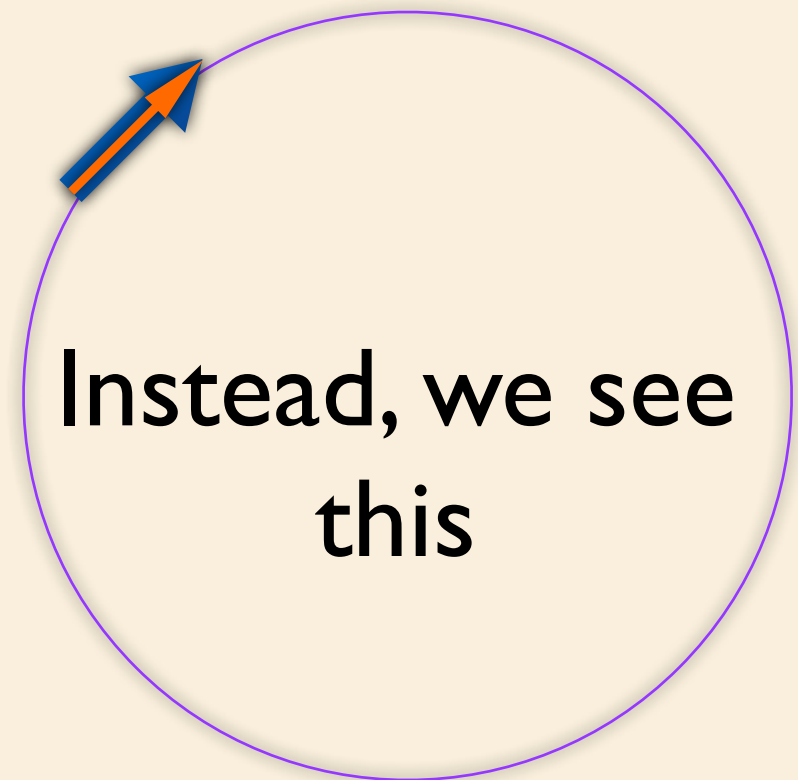
**Muon in a magnetic storage ring**

→ **Muon travel direction**

→ **Muon spin direction**

We don't see this

Instead, we see this
(really this!!)

🔷 **Fermilab**

# The next level of discoveries...

**The next level of discoveries requires the next level of...**



*W* boson 1983



*Top Quark* 1995



*Higgs Boson* 2012

❖ **Fermilab**

# ... the next level of ...

*New physics hides behind particles too heavy to make or too weakly interacting*

interact weakly

Increase intensity

**?New Physics?**

Increase energy

interact strongly

**Known Physics**

light particles

heavy particles

**The next level of Energy**
**The next level of Intensity**
**The next level of Analysis**

**along with these come...**

**The next level of Accelerators**
**The next level of Detectors**
**The next level of Computing**

**The next level of technology**
**The next level of collaboration**

🎗 **Fermilab**

# The **next level** in context of big data



Let's study the evolution of data management and big data

The challenge of Data Management has had many solutions

Examine lessons learned from past and current systems

How can they guide the future of computing, data management and Big Data?

The following is a personal tour.

Allons-y!

**✣ Fermilab**

# The basic problem

Not here, where you need them

Tape Robot

Disk Storage

Compute Farm of worker nodes

**Must identify the files to process**

**Must move those files to the right place**

**Must associate those files with jobs and process them**

**All this must be reliable, scalable, and efficient**

🛠 Fermilab

# My story

HELLO my name is
**Adam**

'02- Staff Scientist at Fermilab in Scientific Computing Division – Head of "Scientific Data Processing" Department, Worked on Tevatron Run 2 D0 and now Muon $g$-2

Was project manager for "SAM" (Data Management system for Fermilab Tevatron Run II experiments and now Intensity Frontier)

'97-02 U Rochester Postdoc on CLEO II.5/III @ Cornell ($b \rightarrow s\gamma$)

'91-97 U Maryland Grad Student on
Tevatron Run 1 D0 @ Fermilab (Supersymmetry)

‡ Fermilab

# The Particle Physics Data Management Story

Why does particle physics push the data management envelope?

The data volume is ever increasing and increasing complexity of detectors

*scalability*

CDF & D0 Run II: ~1M readout channels
Atlas & CMS: ~100M channels

Increasing number of collisions
(bigger haystack to find more needles)
CDF & D0: ~ 10 fb$^{-1}$ (~8 PB) in 10 years
Atlas & CMS: ~25-30 fb$^{-1}$ (~150 PB) in 3 years

🎇 **Fermilab**

# e.g. Evolution in size (note people scale)

Aleph '89–'00



CDF '85–'11



D0 '92–'11



CMS '10–



Atlas '10–

🔷 Fermilab

# Pre-Grid Data Management

**reliability**

o **My experience: Run 1 D0 ('93-97)**

o **FORTRAN, Zebra Banks, HBOOK, PAW**

o **Most processing on laboratory clusters (VAX, SGI, IBM RS-6000, DEC Alpha)**

o **8mm tapes**, mounted by human operators (evolved to small robots, then HPSS) Now use **T10000 tapes @ FNAL**

**efficiency**

o **Some Monte Carlo generation offsite, but not centrally organized**

o **Tapes from remote sites came to Fermilab with traveling students/professors or shipped in mail**

o **Some (lucky) students processed data using sizable resources at their university**

🟢 **Fermilab**

# Pre-grid lessons

o **Don't underestimate the bandwidth of transporting data by truck (or suitcase) [efficient, but not scalable] Modern equivalent is AWS import/export service: (ship your hard drives to Amazon and they'll upload the data to their storage)**

o **Offsite processing was an interesting challenge, often not centrally organized, offsite resources for particular groups**

o **By today's standards not very reliable, scalable, or efficient But good enough to get the job done (discovered top quark)**

| Data Type | DØ Run 1 | CDF Run 1 |
|-----------|----------|-----------|
| Raw | 30 TB, 500 KB/ev | 8 TB, 130 KB/ev |
| Derived | 25 TB | 33 TB |
| TOTAL | 55 TB, 42M events | 41 TB, 64M events |

🎗 Fermilab

# A paradigm shift: transition to C++/OO

My experience (D0 and CLEO '97-02)

FORTRAN did not have the data organization necessary
for the **next level** of data complexity

Industry and education shifting away from FORTRAN

Object oriented programming and C++ were adopted by cadres of
students and postdocs

Lesson Learned: Large scale C++ systems are not easy to write.

Software beginning to become a main part of the experiment, like the
detector hardware
(needs extensive planning, design, construction, upgrades)

An excursion away from the file... OO Databases!

🟢 **Fermilab**

# Objectivity Database

o **Persistency Mechanism:**
   **The form of data on storage (not same as in memory, especially due to pointers)**
   **Translation is overhead – must think hard to get right (efficiency)**

o **Several experiments adopted Objectivity (e.g. BaBar, Cleo III)**
   o **Object Oriented DB Management System**
   o **Automated persistency – retained object structure**
   o **Retrieve events with queries**
   o **Hot (disk) and cold (tape) stores**
   o **A great idea!**
   o **In 2003 BaBar held the world record for the largest civilian hybrid (disk & tape) database – 1 PB**

Fermilab

# But... scalability

For event databases, Objectivity didn't scale:
o Parallel access problematic (hint of things to come)
o $$$
o Large data volume overhead
o Hard to distribute the data
o Federation system difficult
o Reliance on commercial proprietary format

Most experiments returned to the file paradigm by 2004-05
(BaBar – Root; CLEO – PDS [home grown])

Lesson:
A storage system is only as great as how it scales !
An important lesson. Objectivity was an interesting experiment.

🎗 Fermilab

# The rise of ROOT



ROOT
An Object-Oriented Data Analysis Framework

o **Analysis platform and toolkit**
o **Root is to C++ like PAW is to FORTRAN/HBOOK**
o **In '97 I saw first Root demo at Fermilab (reliability)**

o **Embraced data as C++ objects – helped to solve persistency problem**

o **Love it or not, it is ubiquitous in Experimental Particle Physics**

o **My opinion: Dedication of developers to user base was key to success – continues**



o **Root I/O and files are also ubiquitous and have evolved. Having a common format leads to efficiency (e.g. XRootD). Very important to the field.**

o **Lesson: Paradigm shifts require huge commitments.**

**Root survives by always keeping the next level in sight**

**Fermilab**

# Data Handling for Tevatron Run II

**A next level**

**1997 – Joint CDF/DØ Data Management Needs Assessment**

**1999 – Detailed data management plan for D0 (SAM)**

|                  | Plan             | Actual 2012              |              |
|------------------|------------------|--------------------------|--------------|
| Entire dataset   | 0.5 - 1 PB       | 8.2 PB                   | 150x Run I   |
| Tapes            | 50 GB, 6 MB/s    | 800 GB, 120 MB/s         |              |
| Cache            | 20 TB            | 760 TB                   |              |
| # files in catl  | 1M               | 140M                     |              |
| # events in catl | few billion      | abandoned (scale issues) |              |

**Lesson: Think big and prepare for a lot bigger (scalability)**

Fermilab

# Processing



**raw** → Reconstruction

**thumbnail (tmb)**

**tmb** → Correct & Skim

**tmb**

**tmb**
**root-tree** → Make Trees

**root-tree** → Select data

**root-tree**

**root-tree** → Analysis

Detector

Tape storage

**Not shown: Reprocessing, MC generation**

**Plots, Results, Thesis, Paper**

Fermilab

# Behind the scenes

Cache

Compute farm



**Move to disk cache** →

**Move to worker node** →

← **Hold small output files on disk**        **Process data**

**Transfer for merge** →

← **Store big files back to tape**        **Merge small files**

Fermilab

# Behind the scenes

**Tape storage**

**Cache**

**Compute farm**

*Perhaps far away*

*Hopefully not too far away*

**Move to disk cache** →

**Move to worker node** →

← **Hold small output files on disk**    **Process data**

**Transfer for merge** →

← **Store big files back to tape**    **Merge small files**

Data Management Lessons from Particle Physics | Adam Lyon

2014 March 28

**Fermilab**

# Sequential Access via Metadata (SAM)

o **A long evolution: 1998 – Present (16 years!)**
  **Started as a D0 project, then deployed to CDF in '03**

o **In early days storage systems were fragile,**
  **SAM layer provided needed throttling and protection**
  reliability & scalability **– continue to be important lessons**

**Aside: Storage systems can still be fragile even today...**



**Central NFS Storage**

OK

CPUs    Storage Volumes

Bad

CPUs    Storage Volumes

Batch jobs do this!

2014 March 28

**‡ Fermilab**

# SAM for Tevatron Data Handling

○ Forerunner of LHC data management (but important differences on coming slide)

○ *Metadata catalog* (file name, size, events, run info, luminosity, MC details, ...)
   **Lesson:** keep metadata out of the filename

○ Dataset creation and queries
   **Lesson:** need to balance ease of use and flexibility (experts and non-experts)

○ Replica catalog (where all file copies are located)

○ Coordinates and manages data movement to jobs (bbftp, gridftp, dccp, SRM, future XRootD)
   **Lesson:** be ready to adopt new protocols - flexibility is key

○ Cache management (now using dCache)
   **Lesson:** don't do it yourself once something more standard meets requirements

○ File consumption and success tracking, recovery
   **Lesson:** this is an important feature (especially with opportunistic resources)

Data Management Lessons from Particle Physics | Adam Lyon    2014 March 28

🔆 **Fermilab**

# SAM on remote sites (proto-grid)

**circa 2004 - Reprocessing**

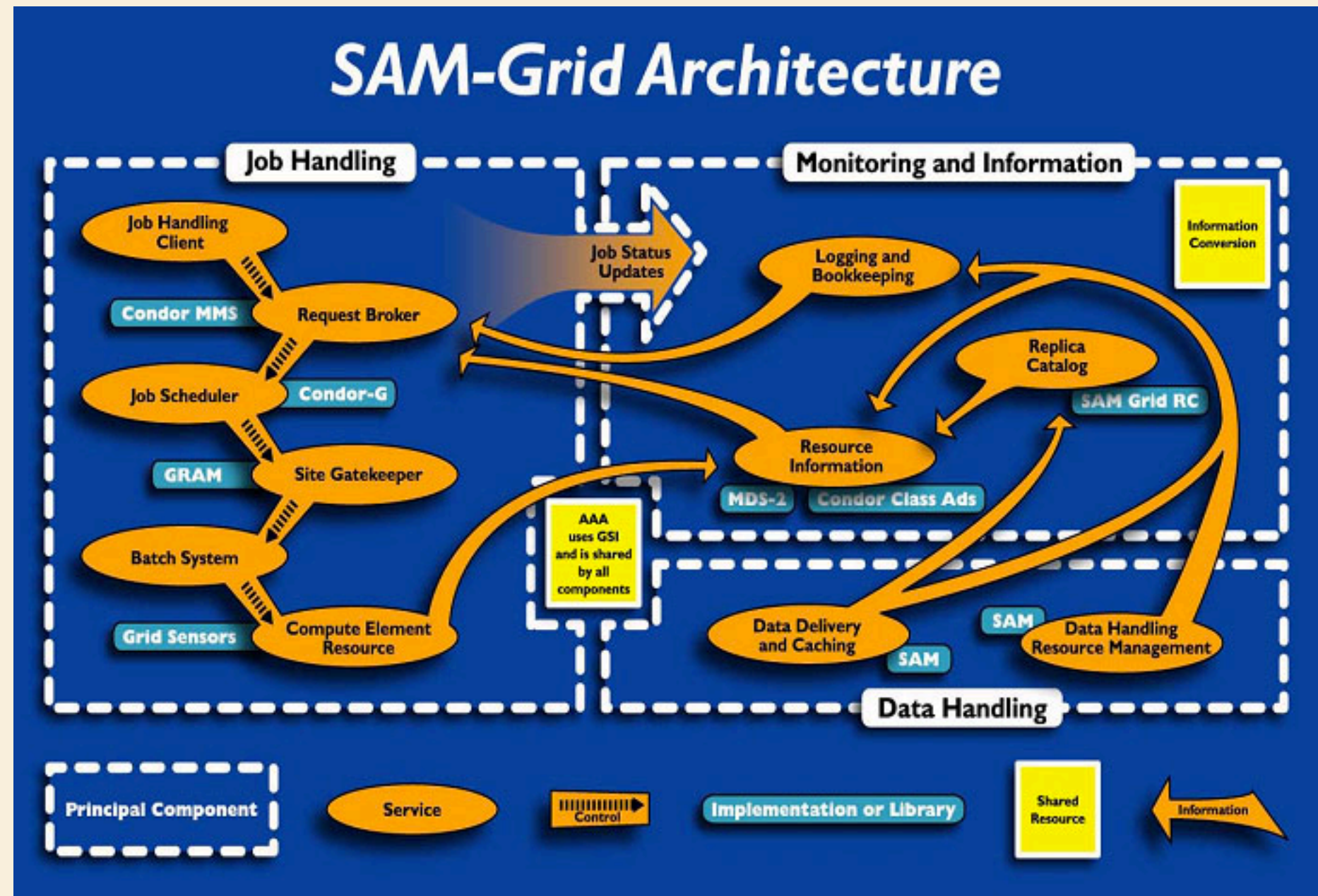**Early sites in US, Europe (Wuppertal, Prague), Brazil**

**Wrote our own middleware while standard middleware was maturing**

**Required a large effort to maintain and operate**

**Lesson: Migrate to standard tools when they mature (much better when someone else maintains and operates). An important theme.**

**Condor, Globus, VDT, OSG**

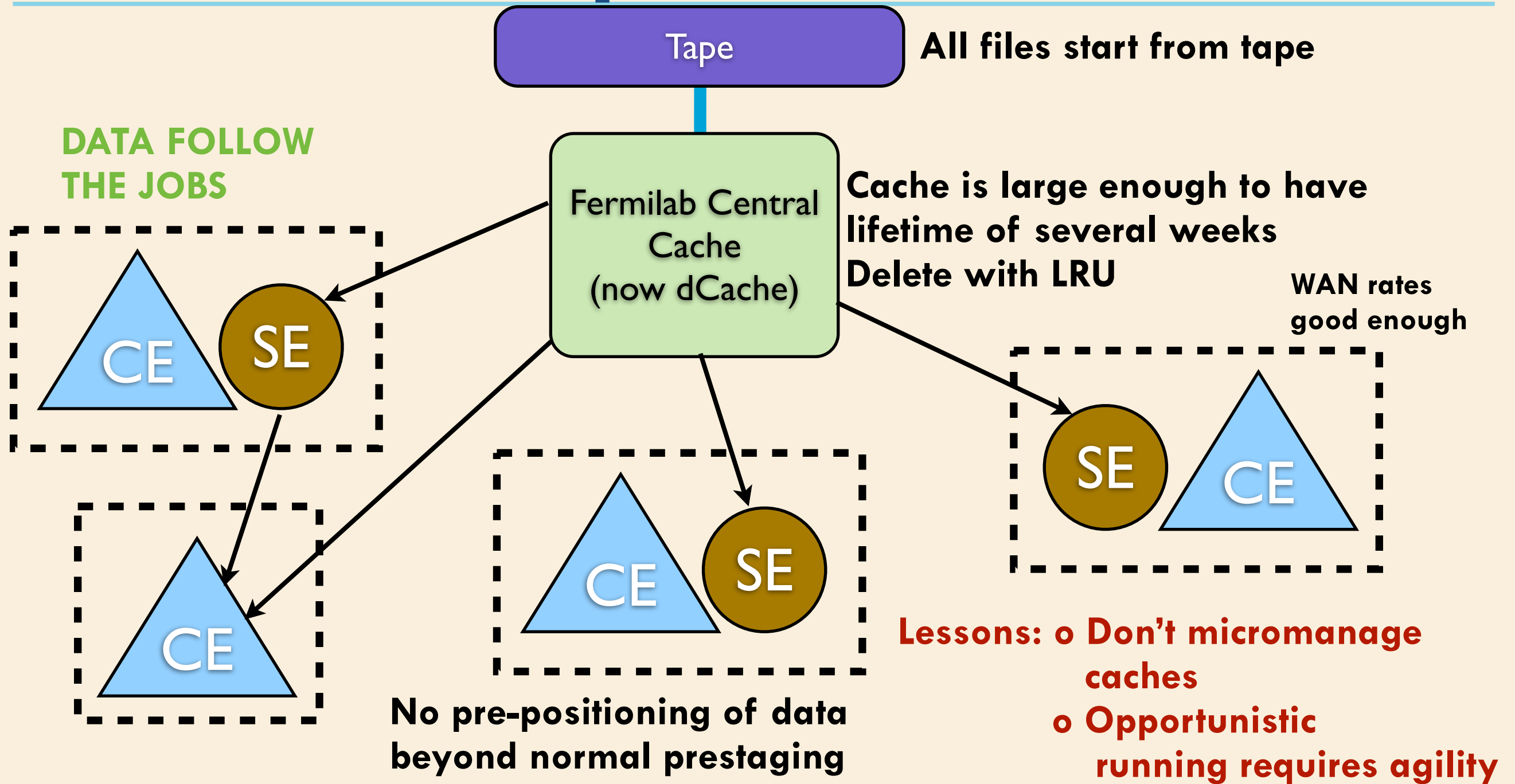**Lesson: Not all sites are the same – fast & slow queues**

# Grid Organizations



o **2007 submitting to OSG/LCG**

o **Support for remote sites is crucial:**
  o **Software stack (condor, globus, vdt)**
  o **Troubleshooting assistance**
  o **Infrastructure**
  o **Security**

o **LHC experiments also provide support**

2014 March 28

🔷 **Fermilab**

# SAM Philosophies

**Tape** — **All files start from tape**

**DATA FOLLOW THE JOBS**

**Fermilab Central Cache (now dCache)**

**Cache is large enough to have lifetime of several weeks Delete with LRU**

**WAN rates good enough**

CE   SE

CE

CE   SE

SE   CE

**No pre-positioning of data beyond normal prestaging**

**Lessons:** o **Don't micromanage caches**
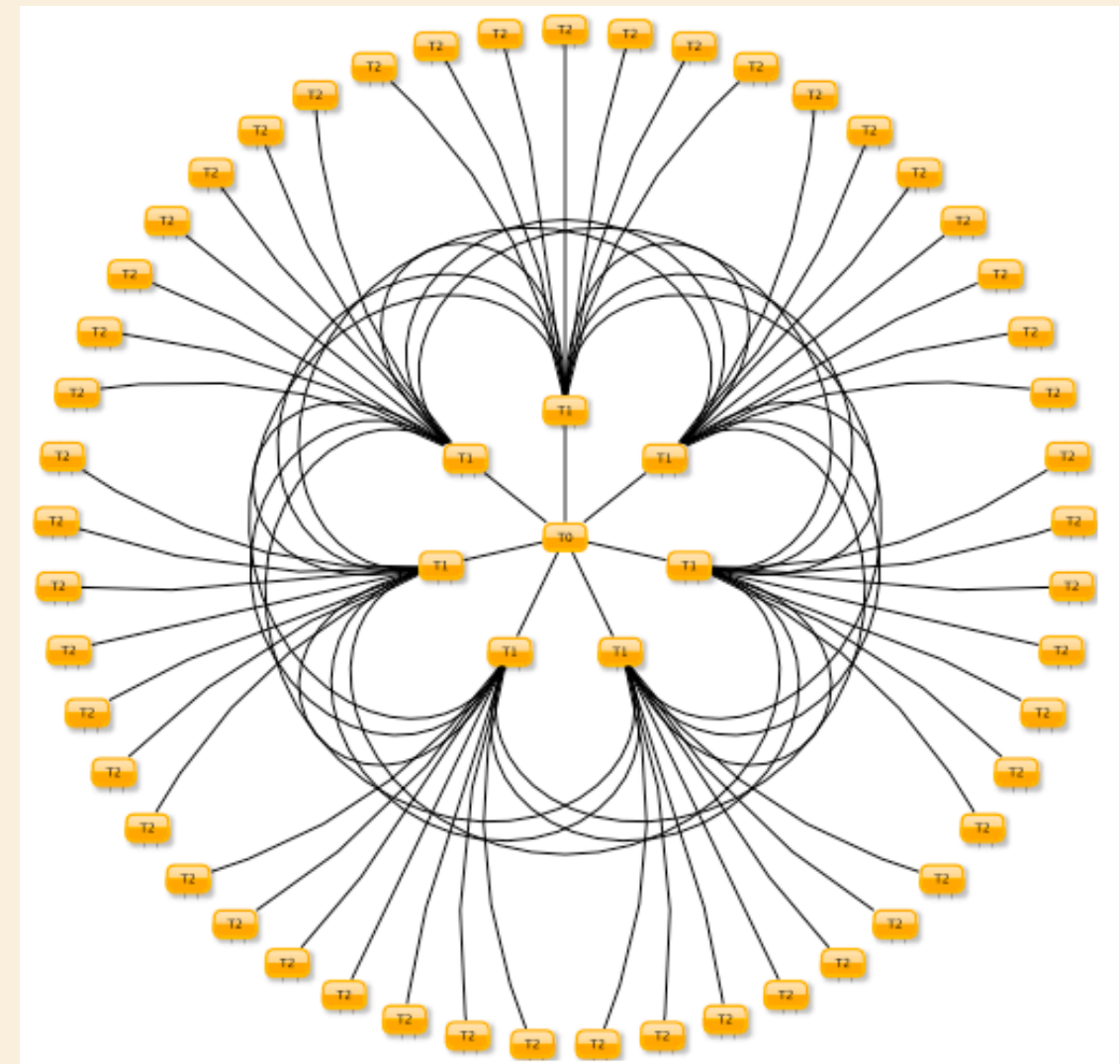o **Opportunistic running requires agility**

## Favor scalability over efficiency

**Run as many sites as we can, even on ones less productive**

**🍀 Fermilab**

# LHC Data Handling (CMS/ATLAS)

o Similar functionality to Tevatron Data Handling. The basic pieces are there, but implemented differently (PhEDEx, DBS3, PANDA, RUCIO, WMAgent/CRAB3, DIRAC)

o LHC Run 1: Hierarchical structure
   o Tier 0 @ CERN: Prompt reco
   o Tier 1 (regional): Re-reco, skimming calibration, make physics format
   o Tier 2 (local): analysis & MC

o Assign datasets to T1's and send jobs to the right place (JOBS FOLLOW THE DATA)

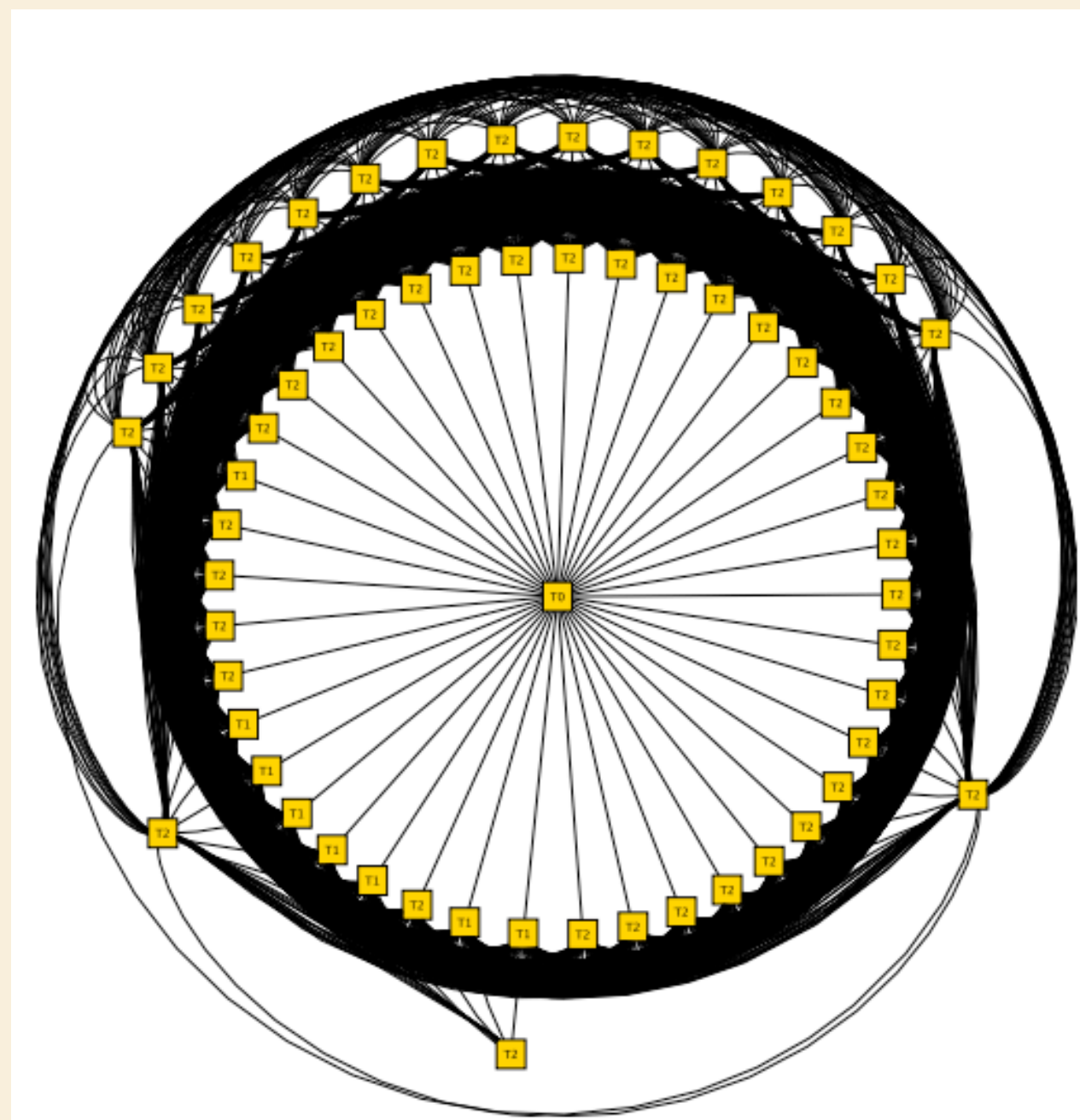o Isolation of tape, limiting data movement allows for slow/non-robust networks, but no opportunistic running

## Favor efficiency over scalability

We control a set of sites, so maximize their productivity

**Fermilab**

# Future LHC Data Handling

o **Preparing for Run 2,** the **next level**
o **Networks are more robust and faster than predicted**
o **Connect everything together**
o **Still pre-position data on T1's but use federated XRootD as a fallback (Any data, Anytime, Anywhere) – <u>mentioned many times here</u>**
o **Makes opportunistic running possible (crucial for LHC Run 2)**



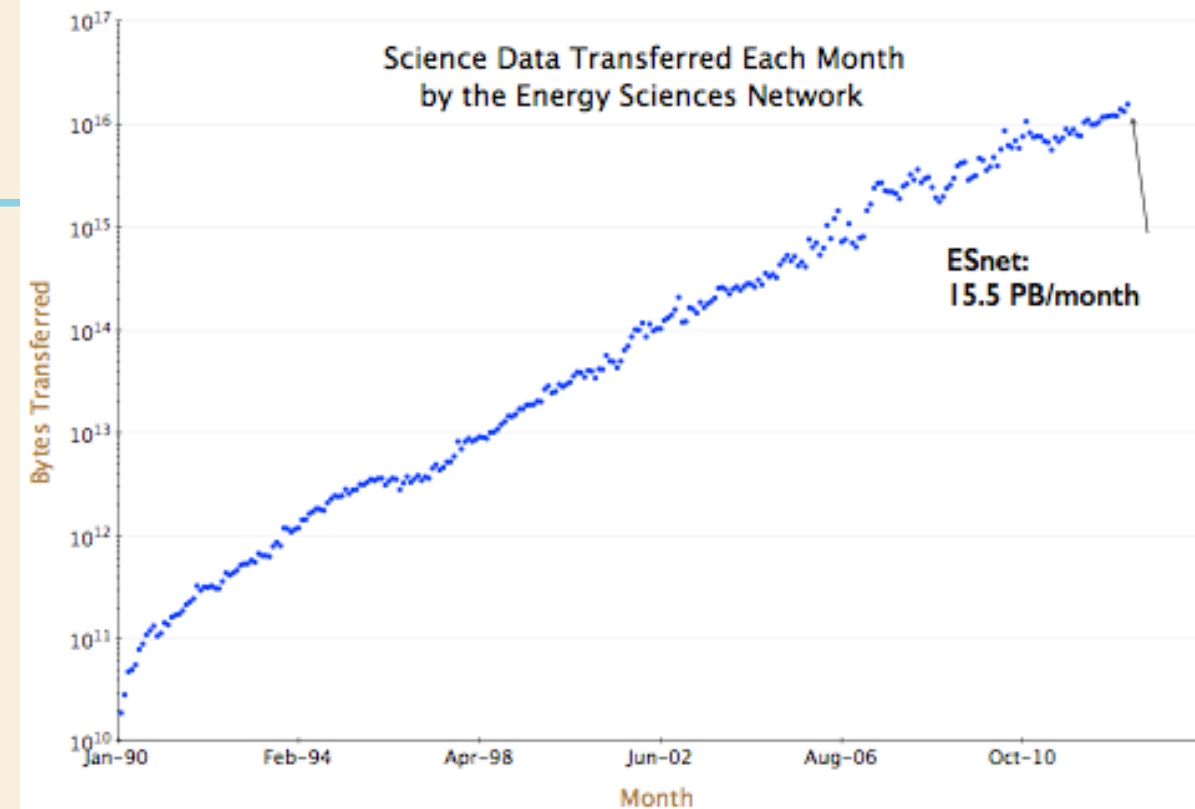### Maximize reliability, scalability, and efficiency
**Wow – the best of all worlds**

**Fermilab**

# Robust Networks

**2003: 10 GB/s trunk lines**
**2014: 100 GB/s trunk lines**

**Large transcontinental pipes**
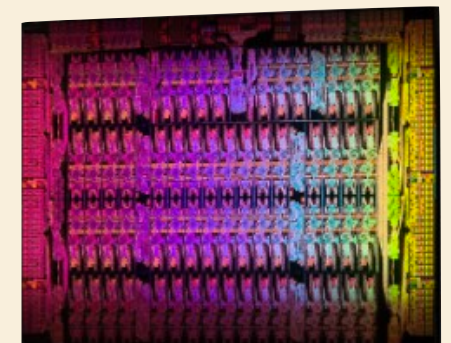
**ESNet From GB to PB**

‡ Fermilab

# Near future: Use new resources



o **Leadership class HPC machines** Make look like an OSG site with OSG-Connect and Parrot

o Take advantage of coming Exascale computing?

o **Cloud resources** to provision worker nodes on demand for peak activity ("Grid Bursting")

But are they suitable for HEP jobs (can be long, high i/o)? Cost effective?

See Garzoglio talk about "On Demand" resources



EC2 Cloud Installation

Fermilab

# The Future

o **The computing landscape is under great change**

o **CPUs can't go faster (no 10 GHz?)**
  - o Limited by power density
  - o So integrate <u>more of them together</u>
  - o Multicore CPUs
  - o GPGPUs (with attached ARMs)
  - o Intel PHI co-processors – Many integrated cores (MICs)

o **Using these technologies is challenging, especially for non-experts**

o **How are data fed to these devices with reliability, scalability, and efficiency?**



Figure 1: Intel CPU Introductions (graph updated August 2009; article text original from December 2004)

# Disruptive technology?

o **Does industry's vision of Big Data (Map-Reduce/Hadoop) have a place in HEP?**

o **Will we need non-traditional High Throughput Computing (HTC) machines (Data Appliances)?**

o **Does the notion of a file survive? E.g. No-SQL databases for events?**
**But remember the lessons from Objectivity!!**

o **Can we use Data Appliances?**
**e.g. EMC Isilon – Very high speed i/o with fast disk caching – can perform map-reduce & DB functions in the cache! Disks become compute nodes! Wow!**

**e.g. YARC – "purpose built" data appliance from Cray**
**Does a "graph analysis" to find connections in data**

o **Must pre-load these appliances. Is opportunistic use possible? How would we use these with reliability, scalability, and efficiency?**

# Is a Paradigm Shift coming?

o **Collider experiments (CMS, Atlas, ...) write event based data. Can this model hold up for the far future? What other models are there?**

o **Does the event model work for the next generation neutrino and muon experiments with trigger-less online? Are they more digital signal processing? Do we need a new set of data management and analysis tools? Just starting...**

o **Collider experiments involve analyzing a subset of the full data (skims). But other experiments may need to analyze EVERYTHING. What do we need to make that work?**

o **What will it take to reach the next level?**

🟦 **Fermilab**

# Summary

o **There's no time for rest!** New experiments & runs are coming and they must succeed!

o Particle physics has an excellent history of grasping the next level of data management to reach the next level of discovery

o Ability and willingness to adjust when reliability, scalability, and efficiency can improve is the key and will allow us to reach that next level

o Are there new machines, technologies, paradigm shifts that take us to the next level with more capabilities than before?

**Fermilab**